

CS 595 - Hot topics in database systems:

## **Data Provenance**

I. Database Provenance

I.1 Provenance Models and Systems

Boris Glavic

November 14, 2012























# Query Languages Overview

## We will cover

- pSQL (DBNotes)
- SQL-PLP (Perm)
- ProQL (Orchestra)
- Color algebra (Mondrian)















# pSQL - the PROPAGATE clause

- Defines how annotations are copied from input to output
- Three different options:
  - Custom
  - DEFAULT
  - DEFAULT-ALL

## Custom

- User defines annotation propagation at the query level
- $R_1.a_1$  TO  $b_1$  = add annotations from input attribute  $R_1.a_1$  to result attribute  $b_1$ .



# pSQL - the PROPAGATE clause

- Defines how annotations are copied from input to output
- Three different options:
  - Custom
  - **DEFAULT**
  - **DEFAULT-ALL**

## DEFAULT-ALL

- Propagate annotations from all cells in the insensitive Where-provenance









# ANNOT semantics

## Sketch Query Evaluation

- Build cross-product of relations and annotation sets in **FROM**
- Apply **WHERE** condition to result
- Evaluate **PROPAGATE** clause for remaining annotations

## Semantics

- If **ANNOT(R.a)** A is only used in **WHERE**
  - $\Rightarrow$  Works as existential check for annotation that fulfill **WHERE**
- If **ANNOT(R.a)** A is also used in **PROPAGATE**
  - $\Rightarrow$  only propagate annotations that fulfill **WHERE**







# pSQL Implementation

- Use relational storage schema for annotated relations
- **Rewrite** pSQL into SQL queries over that schema

# Storage Scheme

- Schema
  - For each attribute  $A$  in relation  $R$
  - Add attribute  $A_a$  that stores annotations on  $A$
- Instance
  - Each tuple can store one annotation on each attribute
  - Duplicate tuples to fit in more annotations

## Example

- $R(a, b)$  will be  $R(a, a_a, b, b_a)$

		<b>R</b>	
		<b>a</b>	<b>b</b>
$r_1$	1	$a_1, a_2$	$2^{a_3}$
$r_2$	2	$2^{a_4}$	3

		<b>R'</b>			
		<b>a</b>	$a_a$	<b>b</b>	$b_a$
$r_1$	1	$a_1$	2	$a_3$	
$r_1$	1	$a_2$	2	-	
$r_1$	2	$a_4$	3	-	



# Translator

- 1 Transform propagate clauses into *Custom* form
- 2 Build bins for each output attribute *b*
  - Add  $R.a$  to bin if  $R.a$  TO  $b$  in propagate clause



# Translator

- 1 Transform propagate clauses into *Custom* form
- 2 Build bins for each output attribute  $b$ 
  - Add  $R.a$  to bin if  $R.a$  **TO**  $b$  in propagate clause
- 3 Generate intermediate queries  $Q_1$  to  $Q_n$ 
  - While at least on bin not empty
  - Take one attribute  $a$  from each bin
  - Generate query that projects each  $a_a$  to  $b_a$
  - Use **NULL TO**  $b_a$  if bin is empty

# Translator

- 1 Transform propagate clauses into *Custom* form
- 2 Build bins for each output attribute  $b$ 
  - Add  $R.a$  to bin if  $R.a$  **TO**  $b$  in propagate clause
- 3 Generate intermediate queries  $Q_1$  to  $Q_n$ 
  - While at least on bin not empty
  - Take one attribute  $a$  from each bin
  - Generate query that projects each  $a_a$  to  $b_a$
  - Use **NULL TO**  $b_a$  if bin is empty
- 4 Generate wrapper query
  - *orderlist* all attribute of pSQL query

```
SELECT DISTINCT *
FROM (Q1 UNION ... UNION Qn)
ORDER BY orderlist
```





# Post-Processor

- Gather all annotations for each result tuple
- Works like aggregation
  - ① Initialize each attribute annotation sets to empty set
  - ② For each tuple add annotations to sets
  - ③ If tuple has different attribute value
    - Output annotated result tuple
    - Start over at (1)









# Example Processing - Post-processing

**result**

	<b>name</b>	<b>name<sub>a</sub></b>	<b>street</b>	<b>street<sub>a</sub></b>	<b>city</b>	<b>city<sub>a</sub></b>
<i>t<sub>1</sub></i>	Bob		10 W 31st	<i>a<sub>1</sub></i>	Chicago	<i>a<sub>1</sub></i>
<i>t<sub>2</sub></i>	Alice	<i>a<sub>2</sub></i>	75 Cermak		Chicago	
<i>t<sub>2</sub></i>	Alice	<i>a<sub>3</sub></i>	75 Cermak		Chicago	
<i>t<sub>3</sub></i>	Gertrud		15 Ellis		Berlin	

# Example Processing - Post-processing

**result**

	<b>name</b>	<b>street</b>	<b>city</b>
$t_1$	Bob	10 W 31st <sup>{<math>a_1</math>}</sup>	Chicago <sup>{<math>a_1</math>}</sup>
$t_2$	Alice <sup>{<math>a_2, a_3</math>}</sup>	75 Cermak	Chicago
$t_3$	Gertrud	15 Ellis	Berlin

**result**

	<b>name</b>	<b>name<sub>a</sub></b>	<b>street</b>	<b>street<sub>a</sub></b>	<b>city</b>	<b>city<sub>a</sub></b>
$t_1$	Bob		10 W 31st	$a_1$	Chicago	$a_1$
$t_2$	Alice	$a_2$	75 Cermak		Chicago	
$t_2$	Alice	$a_3$	75 Cermak		Chicago	
$t_3$	Gertrud		15 Ellis		Berlin	

BY





# Example Translation

## Example

```
SELECT P.name, P.addr
FROM person P,
     ANNOT(P.name) a
WHERE a LIKE '%dumb%'
PROPAGATE P.name TO name
        AND P.addr TO addr
```



```
SELECT DISTINCT *
FROM (SELECT P.name, A.namea, P.addr, P.addra
      FROM person' P,
           (SELECT namea, name, addr FROM person') AS A
      WHERE P.name = A.name AND P.addr = A.addr
           AND A.namea LIKE '%dumb%') AS i
ORDER BY name, addr
```

# Discussion pSQL

- **Extension of existing query language (SQL)**
- For provenance **Generation** and **Retrieval**
- Implemented as **Rewrite**
- Provenance generation always active

# Outline

- 1** Querying Languages for Provenance
  - Querying Provenance Overview
  - DBNotes' pSQL
  - Perm's SQL-PLE
  - ProQL
  - Color Algebra
  - Recap

# SQL-PLE

## Perm Recap

- DBMS with support for provenance generation
- Provenance is generated on-demand
- Implemented as extension of PostgreSQL
- Provenance represented as regular relational data
- Provenance features through **SQL-PLE**



# SQL-PLE Language Overview

## Implementation

- **Rewrite** into SQL

## Language Type

- **Extension of data query language:** super-set of SQL

# Language Features - Provenance Generation

## Request Provenance

```
SELECT prov_expr select_clause  
FROM from_list  
WHERE where_expr  
GROUP BY group_by  
HAVING having_expr  
ORDER BY order_expr
```

- **prov\_expr** = [PROVENANCE [ON CONTRIBUTION (@cs\_type@)]]
  - If present, compute provenance in addition to query results
- **cs\_type** = Which provenance model should be used
  - PI-CS (Perm Influence)
  - C-CS (Perm Copy similar to Where for tuples)
  - Where
  - How (Provenance Polynomials)









# Language Features - Provenance Generation

## External and User-Defined Provenance

- Default: All attributes of relation are used as provenance
- Append **PROVENANCE** ( $A$ ) to **FROM** clause item
- Attributes in list  $A$  are used as provenance
- Use cases
  - External provenance: manually created or from other provenance-aware system
  - Concise representation: Use tuple identifiers instead of complete tuples

# Language Features - Provenance Generation

## Example

```
SELECT PROVENANCE *
FROM (SELECT count(*) AS numCust, country
      FROM customer PROVENANCE (cId) AS c,
           nation PROVENANCE (nId) AS n
      WHERE c.nId = n.nId
      GROUP BY country) AS cCount
WHERE country = 'USA'
```



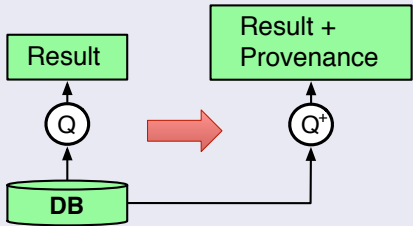
# SQL-PLE Implementation Overview

- **Rewrite**
- Rewrite new language features and provenance generation into plain SQL
- Uses algebraic rewrite rules for provenance generation
- Was discussed in the provenance models class on Perm

# Query Rewrite

## Approach

- Rewrite query  $Q \rightarrow$  query  $Q^+$ 
  - $Q^+$  computes provenance + original results of  $Q$
  - by adding provenance to the inputs (duplicate attributes)
  - propagates provenance through the operations of the query



BY







# Query Rewrite Cont.

## Advantages

- $Q^+$  is single standard SQL
  - $\Rightarrow$  Reuse DB technology
- $Q^+$  accesses the same data as  $Q$ 
  - $\Rightarrow$  No need to store extra information
- Provenance generation as orthogonal SQL extension
  - $\Rightarrow$  Queries over provenance benefit from optimizer
- Correctness formally proven

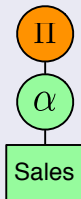


# Example Query Rewrite

Rewrite: 1

## Example Query

```
SELECT total, shop, P(Q+)
FROM
  (SELECT sum(revenue) AS total, shop
   FROM sales GROUP BY shop) AS orig
```



# Example Query Rewrite

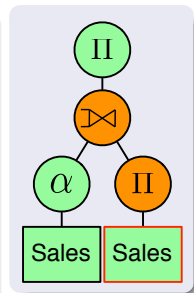
Rewrite: 2

## Example Query

```

SELECT total, shop, P(Q+)
FROM
  (SELECT sum(revenue) AS total, shop
   FROM sales GROUP BY shop) AS orig
LEFT OUTER JOIN
  (SELECT shop AS shop', P(orig+)
   FROM sales+
  ) AS orig+
ON (shop = shop');

```

















# How to Implement Partial Generation

- ( $q$ ) **BASERELATION AS**  $q$  in **FROM** clause
- Apply rewrite rule for base relation to  $q$  instead of rewriting  $q$

## Example

```
SELECT PROVENANCE *
FROM (SELECT count(*) AS numCust, country
      FROM customer c,
           nation n
      WHERE c.nId = n.Id
           GROUP BY country) BASERELATION AS cCount
WHERE country = 'USA'
```







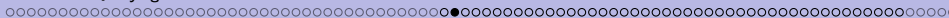


# Outline

## 1 Querying Languages for Provenance

- Querying Provenance Overview
- DBNotes' pSQL
- Perm's SQL-PLE
- ProQL
- Color Algebra
- Recap





# ProQL

## Overview

- Query language for graph provenance representation
- Provenance graphs model How-provenance (provenance polynomials) + mappings (views)
- Graphs are stored in relational database
- Language features for evaluating provenance polynomials in a specific semiring
  - Trust semiring
  - Natural numbers semiring (Bag semantics)
- Main motivation - use in Orchestra



# Schema mappings

- **Schema mapping:** Logical constraints that define the relationship between two schemata
- Different schema may store the same information in different structure
- Schema mappings model these structures in the schema relate
- With some extra mechanism can be use to translate data from one schema into the other
- For now consider mappings as views

## Example

- **Schema**  $S_1$ : `Person(Name, AddrId), Address(Id, City, Street)`
- **Schema**  $S_2$ : `LivesAt(Name, City)`





# Provenance Polynomials Recap

## Rationale

- Use semiring annotations to model provenance
- Annotate a query result tuple with the semiring expression that was used to compute it

## Provenance Polynomials Semiring

- $(\mathbb{N}[I], +, \times, 0, 1)$
- $\mathbb{N}[I]$  - Polynomials with natural number exponents
  - Variables: One per tuple in  $I$
- Convention: annotate each instance tuple with a variable named after its tuple  $id$

BY











# Provenance Polynomials Example II

## Example

$$q = \pi_{Name}(E \bowtie \sigma_{Dep=CS}(P) \bowtie A)$$

$$(q)(t) = \sum_{u.A=t} E(u.E) \times P(u.P) \times (Dep = CS)(u.P) \times A(u.P)$$

Q

Name
Peter
Gertrud

$$e_1 \times a_1 \times p_1$$

$$(e_2 \times a_2 \times p_1) + (e_2 \times a_3 \times p_2)$$

Employee

	Id	Name
$e_1$	1	Peter
$e_2$	2	Gertrud
$e_2$	3	Michael

Assigned

	PName	Id
$a_1$	Server	1
$a_2$	Server	2
$a_3$	Webpage	2
$a_4$	Fire CS	3

Project

	PName	Dep
$p_1$	Server	CS
$p_2$	Webpage	CS
$p_3$	Fire CS	HR

BY

# Provenance in ORCHESTRA

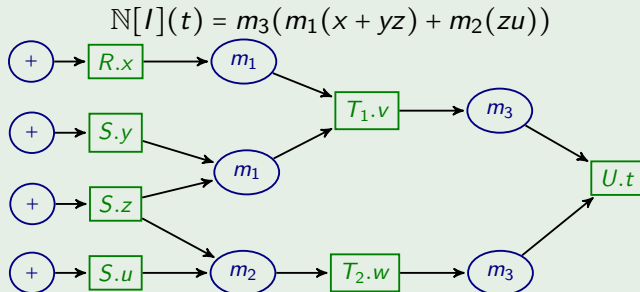
- Use  $\mathbb{N}[I]$  (provenance polynomials)
- Add functions  $m_1, \dots, m_n$  to represent mappings (views)
- E.g.,  $m_3(m_1(x + yz) + m_2(zu))$  means that tuple in view  $m_3$  was derived from
  - tuple in view  $m_1$  derived from  $x, y, z$
  - tuple in view  $m_2$  derived from  $z, u$





## Example Provenance Graph

## Example















# $\leftarrow+\rightarrow$ Expression

- Matches paths of arbitrary length
- Between two tuple nodes

## Example

- $[R] \leftarrow+\rightarrow [S]$ : Matches any path that starts in a tuple node from  $S$  and ends in a tuple node from  $R$
- $\Rightarrow$  paths that connect tuples from  $R$  with their provenance in  $S$

## <ViewName Expression

- Matches a path of length two
- Between two tuple nodes
- Going through view node for ViewName

### Example

- $[\ ] < m_1 [\ ]$ : Matches any application of view  $m_1$
- $\Rightarrow$  paths that connect tuples derived through view  $m_1$  to their provenance tuples from the input of the view

# <Var Expression

- Match paths of length two
- Between two tuple nodes
- Going through any view node
- The view node is bound to Var

## Example

- $[R] \langle \$x \ [] \rangle$ : Matches an application of any view (view node) that uses a tuple from relation  $R$  (incoming edge from tuple node)
- $\Rightarrow$  the view node is bound to  $\$x$











# ProQL Graph Projection Syntax

```
FOR var_binding  
WHERE where_expr  
INCLUDE PATH incl_expr  
RETURN return_expr
```

- **FOR**: bind variables using path expressions
- **WHERE**: logical expression over conditions:
  - Tuple nodes: Conditions on attribute values and relation name
  - View nodes: Test view name
- **INCLUDE PATH**: Paths to copy to output graph
  - Evaluated for each binding of variables that fulfill **WHERE**
- **RETURN**: Variables bindings to return
  - Each combination as tuple

# Examples for Graph Projection Syntax

```
FOR [R $x]  
INCLUDE PATH [$x] <-+ []  
RETURN $x
```

# Examples for Graph Projection Syntax

```
FOR [R $x]  
INCLUDE PATH [$x] <-+ []  
RETURN $x
```

Return the subgraph containing all derivations of tuples in R from base tuples.

# Examples for Graph Projection Syntax

```
FOR [O $x] <-+ [A $y]  
INCLUDE PATH [$x] <-+ [$y]  
RETURN $x
```

# Examples for Graph Projection Syntax

```
FOR [O $x] <-+ [A $y]  
INCLUDE PATH [$x] <-+ [$y]  
RETURN $x
```

Return the part of derivations of tuples in O that involve tuples in relation A.

# Examples for Graph Projection Syntax

```
FOR [$x] <$p [], [$y] <- [$x]
WHERE $p=m1 OR $p=m2
INCLUDE PATH [$y] <- [$x]
RETURN $y
```



# Examples for Graph Projection Syntax

```
FOR [$x] <$p [], [$y] <- [$x]  
WHERE $p= $m_1$  OR $p= $m_2$   
INCLUDE PATH [$y] <- [$x]  
RETURN $y
```

Find tuples that can be derived through mappings  $m_1$  or  $m_2$  and return all one-step derivations from those tuples.

# Examples for Graph Projection Syntax

```
FOR [O $x] <-+ [$z], [C $y] <-+ [$z]  
INCLUDE PATH [$x] <-+ [], [$y] <-+ []  
RETURN $x, $y
```

# Examples for Graph Projection Syntax

```
FOR [O $x] <-+ [$z], [C $y] <-+ [$z]  
INCLUDE PATH [$x] <-+ [], [$y] <-+ []  
RETURN $x, $y
```

Select tuples from O and C that have common provenance and return their derivations.





# Restrictions of Graph Projection

- No creation of new nodes or edges
- No negation
  - This is why we need the + nodes



# Annotation Computation

## EVALUATE semiring OF

- Choose a semiring to evaluate the graph (provenance polynomial over)
  - **Trust**
  - **Bag multiplicity**
  - **Probability**
  - **Other Provenance Models**

## ASSIGN EACH

- Assign semiring annotations to base tuple nodes
  - The ones connected to + view nodes
- Assign functions to views
  - A view function maps lists of annotations to annotations

BY



# Annotation Computation Syntax

```
EVALUATE semiring OF {  
  graph_projection  
}  
ASSIGNING EACH [leaf_node | view_node] variable {  
  CASE case_condition ]: SET annotation  
  [DEFAULT]: SET annotation  
}  
...
```





# Annotation Computation Example

```
ASSIGNING EACH $y {  
  CASE $y in S AND S.a < 3: SET true  
  DEFAULT: SET false  
}
```

# Annotation Computation Example

```
ASSIGNING EACH  $\$p(\$x)$  {  
  CASE  $\$p = m_1$ : SET false  
  DEFAULT  $\$x$   
}
```



# Annotation Computation Example

```
EVALUATE trust OF {
  FOR [0 $x]
    INCLUDE PATH [$x] <-+ [] RETURN $x
}
ASSIGNING EACH leaf_node $y {
  CASE $y in C : SET true
  CASE $y in A AND $y.height >= 6 : SET false
  DEFAULT : SET true
}
ASSIGNING EACH mapping $p($z) {
  CASE $p = m4 : SET false
  DEFAULT : SET $z
}
```

# Implementation

- Use relational representation of provenance graph
  - Reuse existing relations for tuple nodes
  - Add new relations for edges and view nodes
- Translate ProQL into SQL (datalog) queries





# Translation of ProQL into SQL

- 1 Build a **provenance schema graph** based on views
  - Independent of query
- 2 Match query to schema graph
  - Which nodes in schema graph may match path expressions
- 3 Create datalog program based on matching
- 4 Execute program using DBMS

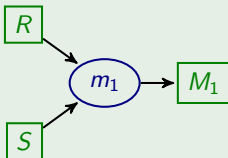
# 1) Provenance Schema Graph

- Models possible provenance relationships among tuples
- One node for each view
- One node for each relation
- Directed edge from a relation to view
  - View accesses this relation
- Directed edge from view to relation
  - View generates this relation
- **Property:** Each path in the a provenance graph corresponds to path in schema graph
  - by mapping each tuple node to its relation
  - by mapping a view application node to the view node

# 1) Provenance Schema Graph

## Example

```
CREATE VIEW  $M_1$   
SELECT b AS e FROM R JOIN S ON (R.a = S.c)
```





# Example Schema Graph

## Example

```
CREATE VIEW T1 AS (  
SELECT * FROM R  
UNION  
SELECT a.* FROM S a, S b WHERE a.x = b.y)
```

```
CREATE VIEW T2 AS (  
SELECT a.* FROM S a, S b WHERE a.y = b.y)
```

```
CREATE VIEW U AS (  
SELECT * FROM T1  
UNION  
SELECT * FROM T2)
```



## 2) Matching Path Expressions

- For each path expression in the query
- Start with all relation nodes that match the start point
- Compute all potential paths in the schema graph that match the path expression
  - Restriction on views
  - Restriction on end points





### 3) Generate Datalog Program

- Build rules for each path in schema graph for each variable
  - Edges are tuples from provenance tables
  - Joins between provenance tables are paths
- Integrate `WHERE` conditions
- Result tables
  - **Provenance tables**: store the edges in the result graph
  - **RETURN bindings**: store variable bindings from `RETURN` clause





# Example Generate Datalog Program

## Example

```
FOR [$z] <- [$y], [$y] < $m_1$  [$x]
INCLUDE PATH [$z] <-+ [$x]
RETURN $x, $z
```

### Provenance Tables

Two tables for  $m_1$ :  $m_1(r, t_1)$  and  $m'_1(s_1, s_2, t_1)$

One table for  $m_2$ :  $m_2(s_1, s_2, t_2)$

Two tables for  $m_3$ :  $m_3(t_1, u)$  and  $m'_3(t_2, u)$

### Datalog Rules for RESULT

$$\text{return}(x, z) : -m_1(x, y) \wedge m_3(y, z)$$
$$\text{return}(x, z) : -m'_1(x, u, y) \wedge m_3(y, z)$$
$$\text{return}(x, z) : -m'_1(u, x, y) \wedge m_3(y, z)$$



## 4) Execute over DBMS

- Add **CASE** constructs for ASSIGNING EACH clause
- Determine how provenance polynomial structure using paths
- **UNION ALL** instances of paths
- Use aggregation functions to combine annotations
- If provenance graphs do not have cycles
  - Unfold the recursive datalog rules
  - ⇒ Eliminates recursion

# Example

## Example

### Datalog Rules for RESULT

*return*( $x, z$ ) :  $\neg m_1(x, y) \wedge m_3(y, z)$

*return*( $x, z$ ) :  $\neg m'_1(x, u, y) \wedge m_3(y, z)$

*return*( $x, z$ ) :  $\neg m'_1(u, x, y) \wedge m_3(y, z)$

```
SELECT r AS x, u AS z FROM m1, m3 WHERE m1.t1 = m3.t1
UNION ALL
SELECT s1 AS x, u AS z FROM m'1, m3 WHERE m'1.t1 = m3.t1
UNION ALL
SELECT s2 AS x, u AS z FROM m'1, m3 WHERE m'1.t1 = m3.t1
```





# Outline

## 1 Querying Languages for Provenance

- Querying Provenance Overview
- DBNotes' pSQL
- Perm's SQL-PLE
- ProQL
- Color Algebra
- Recap



# Color Algebra

## Overview

- Query language for annotation propagation (colors)
- Data model: Relational with annotations
  - Annotations on **blocks** of values spanning multiple attributes
- Language is an extension of relational algebra
  - Standard operators that propagate annotations
  - Operators that access annotations
- Annotation propagation always active







# Projection

- Projection  $\pi_A(R)$ 
  - $A$  is a set of attributes
- Standard relational projection for data
- Keep all annotations on blocks overlapping with  $A$

## Definition (Projection)

$$\pi_{A_1, \dots, A_k}(R)$$

$$R' = \{t[A_1, \dots, A_k] \mid t \in R\}$$

$$\chi'(t[A_1, \dots, A_k], Y) = \bigcup_{Z \in R} \chi(t, Y \cup Z)$$

BY

## Projection

## Example

 $\pi_{Name}(Employee)$ 

## Employee

	Name	Salary
$e_1$	Peter	\$30
$e_2$	Helga	\$50
$e_3$	Ann	\$199

$$\chi(e_1, \{Name\}) = \{Manager\}$$

$$\chi(e_1, \{Salary\}) = \{CAD\}$$

$$\chi(e_3, \{Name, Salary\}) = \{Since 2000\}$$

## Q

	Name
$t_1$	Peter
$t_2$	Helga
$t_3$	Ann

$$\chi(e_1, \{Name\}) = \{Manager\}$$

$$\chi(e_3, \{Name\}) = \{Since 2000\}$$



# Selection

- Selection  $\sigma_C(R)$
- $C$  is an equality comparison like
  - 1  $A = a$  where  $A$  is attribute and  $a$  a constant
  - 2  $A = B$  where  $A$  and  $B$  are attributes
- Standard relational selection for data
- For (1) Keep all annotations on tuples  $t \models C$
- For (2) Only keep colors  $c$  for exists blocks
  - Block containing  $A$  of color  $c$
  - Block containing  $B$  of color  $c$





## Selection

## Example

 $\sigma_{Name=Nickname}(Names)$ 

Names

	Name	Nickname
$n_1$	Peter	Fatty
$n_2$	Helga	Rose
$n_3$	Ann	Ann

 $\chi(n_1, \{Name\}) = \{Middle\ name\}$  $\chi(n_1, \{Nickname\}) = \{Why?\}$  $\chi(n_3, \{Name, Nickname\}) = \{Boring\}$ 

Q

	Name	Nickname
$t_1$	Ann	Ann

 $\chi(t_1, \{Name\}) = \{Boring\}$  $\chi(t_1, \{Nickname\}) = \{Boring\}$







# Cross Product

## Example

$Employee \times Names$

		Q			
		Name	Salary	FirstName	Nickname
$t_1$		Peter	\$30	Ann	Ann
$t_2$		Helga	\$50	Ann	Ann
$t_3$		Ann	\$199	Ann	Ann
		...	...	...	...

$$\chi(t_1, \{FirstName\}) = \{Manager\}$$

$$\chi(t_1, \{Name, Nickname\}) = \{Boring\}$$

$$\chi(t_2, \{Name, Nickname\}) = \{Boring\}$$

$$\chi(t_3, \{Name, Nickname\}) = \{Boring\}$$



## Block Projection (Lower)

- Lower Block Projection  $\pi_A^L(R)$ 
  - $A$  is a set of attributes
- Keep only tuples that have at least one annotation on a superset of  $A$
- Keep only annotations on blocks covering  $A$

### Definition (Block Projection (Lower))

$$\pi_A^L(R)$$

$$R' = \{t \mid t \in R \wedge \exists Y : A \subseteq Y \wedge \chi(t, Y) \neq \emptyset\}$$
$$\chi'(t, Y) = \begin{cases} \chi(t, Y) & \text{if } A \subseteq Y \\ \emptyset & \text{otherwise} \end{cases}$$

BY



# Block Projection (Upper)

- Upper Block Projection  $\pi_A^U(R)$ 
  - $A$  is a set of attributes
- Keep all tuples
- Keep only annotations on blocks contained in  $A$

## Definition (Block Projection (Upper))

$$\pi_A^L(R)$$

$$R' = R$$

$$\chi'(t, Y) = \begin{cases} \chi(t, Y) & \text{if } Y \subseteq A \\ \emptyset & \text{otherwise} \end{cases}$$

BY







# Block Selection

## Example

 $\sigma_{Manager}(Names)$ 
**Employee**

	Name	Salary
$e_1$	Peter	\$30
$e_2$	Helga	\$50
$e_3$	Ann	\$199

$$\chi(e_1, \{Name\}) = \{Manager\}$$

$$\chi(e_1, \{Salary\}) = \{CAD\}$$

$$\chi(e_3, \{Name, Salary\}) = \{Since\ 2000\}$$

**Q**

	Name	Salary
$t_1$	Peter	\$30

$$\chi(t_1, \{Name\}) = \{Manager\}$$

## Merge

- Merge  $\mu_{X,Y}(R)$
- $X$  and  $Y$  are sets of attributes from  $R$
- Keep all tuples
- Merge annotations on blocks contained in  $X$  and  $Y$

## Definition (Merge)

$$\mu_{X,Y}(R)$$

$$R' = R$$

$$\chi'(t, Z) = \{c \mid \exists X_1, X_2 : Z = X_1 \cup X_2 \wedge X_1 \subseteq X \wedge X_2 \subseteq Y \\ \wedge c \in \chi(t, X_1) \wedge c \in \chi(t, X_2)\}$$

## Merge

## Example

 $\mu_{\{Name\}, \{Salary\}}(Employee)$ 

Employee

	Name	Salary
$e_1$	Peter	\$30
$e_2$	Helga	\$50
$e_3$	Ann	\$199

 $\chi(e_1, \{Name\}) = \{Manager\}$ 
 $\chi(e_1, \{Salary\}) = \{CAD\}$ 
 $\chi(e_3, \{Name\}) = \{Since\ 2000\}$ 
 $\chi(e_3, \{Salary\}) = \{Since\ 2000\}$ 

Employee

	Name	Salary
$e_1$	Peter	\$30
$e_2$	Helga	\$50
$e_3$	Ann	\$199

 $\chi(t_3, \{Name, Salary\}) = \{Since\ 2000\}$



# Union

- Union  $R \cup S$  with color functions  $\chi_R$  and  $\chi_S$
- Apply standard relational union to data
- Keep all annotations from both inputs

## Definition (Union)

 $R \cup S$ 

$$R' = R \cup S$$

$$\chi'(t, Y) = \chi_R(t, Y) \cup \chi_S(t, Y)$$



# Implementation in Mondrian

## Relational Representation

- Extend each relation  $R$  with additional attributes to store annotations
  - A boolean attribute  $bA$  for each attribute  $A$  of  $R$  to store block memberships
  - One attribute  $\lambda$  to store the actual annotation
- Duplicate tuples if necessary (Set semantics)

## Query rewrite

- Rewrite color algebra expressions into SQL

# Example Relational Representation

## Example

		Employee	
		Name	Salary
$e_1$		Peter	\$30
$e_2$		Helga	\$50
$e_3$		Ann	\$199

$$\chi(e_1, \{Name\}) = \{Manager\}$$

$$\chi(e_1, \{Salary\}) = \{CAD\}$$

$$\chi(e_3, \{Name, Salary\}) = \{Since\ 2000\}$$

## rep(Employee)

		Name	Salary	bName	bSalary	$\lambda$
$e_1$		Peter	\$30	1	0	Manager
$e_1$		Peter	\$30	0	1	CAD
$e_2$		Helga	\$50	0	0	-
$e_3$		Ann	\$199	1	1	Since 2000









# Outline

- 1** Querying Languages for Provenance
  - Querying Provenance Overview
  - DBNotes' pSQL
  - Perm's SQL-PLE
  - ProQL
  - Color Algebra
  - Recap



# Recap

## Motivation for Querying Provenance

- Size of Provenance Information
- Extract parts of interest from provenance

## Challenges of Querying Provenance

- Different model for data and provenance
- Queries that span regular data and provenance
- Different access patterns



# Recap

## Language Design

- Support Retrieval and/or Generation features
- Extending the data query language
- ... or Develop new language

## Language Implementation

- Compile into data query language
- Create new or modify execution engine



# Recap

## Covered Languages

- **pSQL**
  - Generation and retrieval language for annotations
  - Extension of subset of SQL compiled into SQL
  - Annotation propagation always active
  - **DBNotes**
- **SQL-PLE**
  - Generation language for provenance
  - Retrieval using SQL
  - Extension of SQL
  - **Perm**
- **ProQL**
  - Retrieval language for provenance graphs ( $\mathbb{N}[I]$ )
  - New query language compiled into SQL
  - **Orchestra**

BY



# Recap

## Covered Languages

- **Color Algebra**
  - Generation and retrieval language for annotations
  - Extension of relational algebra compiled into SQL
  - Annotation propagation always active
  - **Mondrian**





# Literature II



[Boris Glavic.](#)

Perm: Efficient Provenance Support for Relational Databases.  
PhD thesis, University of Zurich, 2010.



[Laura Chiticariu, Wang-Chiew Tan, and Gaurav Vijayvargiya.](#)

DBNotes: a Post-it System for Relational Databases based on Provenance.  
In SIGMOD '05: Proceedings of the 31th SIGMOD International Conference on Management of Data, 942–944, 2005.



[Deepavali Bhagwat, Laura Chiticariu, Wang-Chiew Tan, and Gaurav Vijayvargiya.](#)

An Annotation Management System for Relational Databases.  
The VLDB Journal, 14(4):373–396, 2005.



[Floris Geerts, Anastasios Kementsietsidis, and Diego Milano.](#)

MONDRIAN: Annotating and Querying Databases through Colors and Blocks.  
Technical report, University of Edinburgh, 2005.



[Floris Geerts, Anastasios Kementsietsidis, and Diego Milano.](#)

iMONDRIAN: A Visual Tool To Annotate and Query Scientific Databases.  
In EDBT '06: Proceedings of the 9th International Conference on Extending Database Technology (demonstration), 1168–1171, 2006.